

Plot Based Movie Recommendation System

Dr. Jyothi R

Associate Professor

Department of Computer Science and Engineering

PES University Bangalore, India

jyothir@pes.edu

Ipshita Biswas

Department of Computer Science and Engineering

PES University Bangalore, India

ipshitabiswas117@gmail.com

Sunil Parameswaran

Department of Computer Science and Engineering

PES University Bangalore, India

sunil05.parameswaran@gmail.com

Anjan R Prasad

Department of Computer Science and Engineering

PES University Bangalore, India

anjanprasad.21@gmail.com

Prarthana Jyothi

Department of Computer Science and Engineering

PES University Bangalore, India

jprats302@gmail.com

Abstract—This project is a giant leap in the field of movie recommendation systems, a field that has been really stubborn when it comes to getting good results out of estimating user preferences. The current platforms usually end up too generic with its recommendation feature, causing the user only frustration and dropping engagement. With useful insights into the movie plots coupled with the use of emerging technologies like Word2Vec for semantic analysis, vectorization for data representation, R2N2 for sequential processing, batch normalization to bring model stability, and NLP for overall text understanding, this system provides a lot more personalized and precise movie recommendations. What is more, it goes beyond the mere recommendation by mastering the categorization of genre through analyzing numerous movies, hence ensuring users discover related movies following their taste and probably diversify into other genres kinds of movies that are to their preference. This site is ultimately here for one sole purpose: to give an enhanced overall experience in the discovery of movies, so that end users may easily and with satisfaction find and enjoy some of their favorite films.

Index Terms—Word2Vec, Vectorization, RNN, Sequential, Batch Normalization, NLP

I. INTRODUCTION

This is not just a project; it is more of an in-depth study at the fascinating crossroad of technology and entertainment, with special focus pertaining to the emergent field of movie recommendation systems. This is because, in the rapidly changing current digital environment, streaming services are increasing rapidly in an extraordinary way, and the future demand for intelligent algorithms will be much higher. These algorithms, in essence, will look at very large libraries of different kinds of content to give a suggested movie recommendation that best matches the preference and liking of a given user to very particular detail.

Importance: The importance of this project is in its ability to adapt and thereby meet the diversified and wide-ranging user-tastes and preferences. In essence, it aims to offer a much more engaging user experience, while at the same time offering

fine curation of films that adhere very closely to the unique interests of the user. It can humanly understand and learn user preferences in a much better way than current systems that do not catch the drift of what users want.

The backbone of the project: the most sophisticated natural language processing (NLP) algorithms running on top of powerful deep neural networks. Working in tandem, these technologies dissect and process with scrupulous precision each intricate storyline that has been featured in the movies. The process is defined as the transformation of narrative content into structured numerical representations through the use of sophisticated methodologies such as TF-IDF, stemming, and vectorization. The Deep Neural Network then processes these numerical forms to reveal intricate relationships and themes which may not, at the face of it, be immediately obvious even to a seasoned analyst.

This project aims to push the envelope of what the existing recommendation systems can provide for the experience of movie-watching: deeply personalized and, by this very fact, profoundly immersive, to set up new standards in the industry. Empower the users to go beyond their regular cinematic likings and enable them to discover some genres or films which will really reflect their personal taste and culture interest.

Our vision is to try to contribute to creating an environment that will do much more than simply help in finding interesting content but build a bridge between users and the great, immense world of cinema. The platform will encourage the deepest relation between film lovers and the films they love by the largest motivation of discovery ever. Through this project, we hope to redefine how people are able to interact with film and make it a more personalized, inclusive, and rewarding piece of their digital lives.

II. LITERATURE SURVEY

The revolutionary change from simple collaborative and content-based filtering systems at the beginning of the movie recommendation system to the advanced integration of several machine learning and deep learning techniques reflects the game-changer toward improving user experience in the consumption of digital content. This meant the systems, crucial in personalizing users' interactions with the streaming platforms, had migrated from analyzing explicit user feedback: things like ratings, for example, to discovering complex behavior patterns on the part of the user. This consolidation of research examines methodologies, algorithms, and evaluation metrics [2] that define present recommendation systems, among others, to underscore the position of such systems in navigation of the large and expanding digital content landscape and the continued search for higher accuracy and user satisfaction.

A. Evolution of Recommendation Systems

This is where the journey of recommendation systems starts, featuring collaborative and content-based filtering, mainly via simple user-item interactions. In this system, early ones used explicit user feedback, such as rating, for predicting the user's preferences. With the advent of machine learning and deep learning, systems have largely improved; they can figure out complex patterns among many and other latent factors affecting user preferences.

B. Methodologies in Movie Recommendation

Collaborative Filtering (CF): Traditional CF techniques, such as User-Item Matrix Factorization, have laid a stronger foundation for recommendation systems. The basis of this approach is on the similarity between users or items for the purpose of developing recommendations.[10] Content-Based Filtering: This filtering approach recommends items of the same type that a user liked in the past, based on item features, such as the genres of movies or their plots.[3] Hybrid Systems: Combining CF with content-based methods, hybrid systems strengthen the reasoning of each method to get more accurate and diversified recommendations.

C. Advanced Machine Learning Techniques

More recent recommendations in the recommendation systems use cutting-edge algorithms, particularly deep learning and natural language processing (NLP).[3] Conventional models, like Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), have been harnessed to capture sequential and temporal dynamics within user interactions. [9] On the other hand, NLP is of the essence when analyzing movie plot summaries and reviews to bring out a proper perspective on what the user sentiment and preference are. Such improvements make the approaches more effective in content-based filtering with regard to the recommendation system.

D. Evaluation Metrics

The effectiveness of recommendation systems is measured through various metrics: Some of the commonly used metrics to evaluate recommendations are Mean Squared Error (MSE) and Root Mean Squared Error (RMSE) [4]. Diversity and novelty are metrics of recommendations, which allow accuracy at the same time that they allow representing a larger space of the variety of user interests. So, in the end, its success lies at the level of user engagement and satisfaction, which is usually measured via A/B testing and user surveys, respectively.

E. Case Studies and Applications

Some previous studies even found that the application of these methods, like those reviewed in this survey, is validated through real-world implementation.[11] For example, some keyword- and plot-based content systems using keywords already have a way of efficiently recommending similar movies, while user feedback will improve the algorithms.

F. Challenges and Future Directions

There is, of course, a lot of work to do with solving the common challenges of recommender systems, from data sparsity to the cold-start problem and providing a balance between accuracy and diversity. If nothing else, future research has to look forward to integrating much more user-centric data and enhancing the models with cross-domain knowledge; their realms should be pried open for privacy issues in personalized recommendations.

From very basic algorithms, the movie recommendation systems field has shifted to the latest and most sophisticated machine learning models that promise much improvement to user experience. Continuous research and development in this field does, however, promise much more personalized and accurate recommendation systems in the future.

III. OBJECTIVES

- Design and develop a content-based filtering model using movie plot summaries for generating recommendations. The objective of the model is to identify and recommend films users would like to watch on the same theme or storyline provided as input to them.
- To Improve User Experience through Personalization: Improve user experience for the digital streaming platforms by providing users with movie recommendations tailored for individual users.[12] The system shall try to learn and adapt to the individual users' preferences for the system, and therefore, improve users' engagement and satisfaction.
- Use Natural Language Processing (NLP) to harness the power of state-of-the-art NLP techniques in the analysis and extraction of meaningful features from the summaries of movie plots. Textual data are processed to bring out elements that will make the movies similar to each other.
- Increase Recommendation Accuracy: The main goal is to increase recommendation accuracy by further fine-tuning recommendation algorithms that are used in determining

similarity between plots of different movies. Relevance and high affinity for user taste should be reached close to this goal with recommendations.

- Addressing Data Sparsity with Diversity: To better deal with the problem of data sparsity when it comes to recommendations, the site opts to go with a blanket effect. From what is most popular to perhaps less popular, niche content, but for the user to feel he can find anything within our service.
- System performance assessment: The recommendation system performance is tested against how effectively it gives outputs to proper system performance criteria, that is, precision, recall, and F1 score. The system will be adapted based on system performance assessment for improvement and received feedback.
- Advanced Machine Learning Models: Analyze how advanced machine learning models, including deep learning and neural networks, can be applied to further the recommendation system for bigger datasets' capability and making better predictions.

IV. DATASET

A. Features

The TMDB 5000 movies dataset is encompassed in the information that covers an approximate number of 5000 movies, including their title, genres, budget, revenue, release date, and rating. This gives an amount of information that details out the commercial, as well as artistic sides of the film industry.

Of these attributes in the movie world, for example, the dataset includes the title name for each movie. For the genre of the world of film, the categorization of themes or style might be a basic characteristic or identifier; for example, one could have action or comedy. Keywords related to every movie reveal the core themes and narrative constituents that contribute to the understanding of the content of the film. Financial details like budget, revenue, and release date, reveal insight into the economic aspects and lifecycle of a movie in the market. Ratings gauge audience reception and popularity, reflecting the film's success and acceptance.

B. Data Cleaning

But data cleaning from the TMDB dataset is systematically done that allows making the dataset fit for the process of analysis that follows. First of all, the data has been parsed with the main intention of reducing the complexity of the data structures. This particularly streamlined the columns like "genres" and "keywords". These were custom transformation functions that operated them to a form managing which represents genres and keywords as lists of dictionaries with stripped strings and commas delineating individual elements. This transformation made for a brief informative summary with reference to the movie, which gives an idea of the genres and keywords of the movie. Now it was possible to dig deep into what was represented.

We have then excluded all columns that are not required and will have no bearing on the dataset: 'homepage', 'original_title', 'status', 'spoken_languages'. Critique will be made on these irrelevant columns with respect to relevance and fit in the scope and magnitude of the analysis. There are those which are probably not needed or extraneous to the analytical objectives. These were removed as an attempt to make an effort to clean up the dataset so as the resultant would be analyzable and manageable.

The essential step of data cleaning included the need to impute any missing or null values so that the data is whole and complete.[15] The exact methods, thus, were not clearly defined with respect to handling missing values in the data. However, it can be inferred that the author, in this study, must have applied a plethora of methods—from deletion and imputation to further exploration regarding possible sources which could have caused the data gaps.

This was done in view of ensuring that the cleaning steps that follow would not compromise the accurate and representative nature of the original data. Validation measures were quite strenuous to maintain the reliability and strength of the data; however, there is high confidence in the suitability of the data for the tasks in subsequent analysis. Generally, the data cleaning process in TMDB dataset took care of the details in the transformation, pruning, and validation of the data. It would be a result with a very solid foundation for meaningful analysis and the understanding of trends and patterns within the underlying TMDB data, assuring completeness, integrity, and relevance of data.

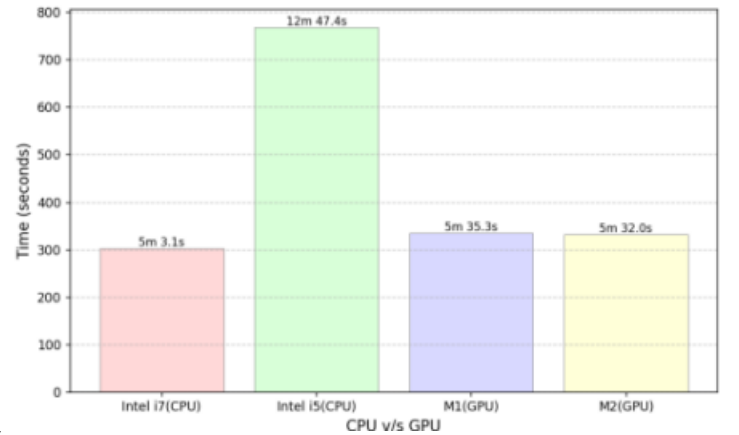


Fig. 1. Training time on CPU vs GPU

V. IMPLEMENTATION

The carefully outlined implementation procedure in this paper constitutes a complete workflow for preprocessing textual data, training a Word2Vec model, and building a neural network model for genre prediction. This provides efficient, accurate, and robust designs at each stage, including data preprocessing, model training, and inference.

From the TensorFlow Device Configuration step, the implementation looks through the available computational resources

to see to it that the GPU is present for much-accelerated improved training. This, in effect, guarantees effective utilization of the hardware resources and greatly influences the efficiency of computations and, therefore, training time. This is effective for model training and efficient inference implementation since it uses available hardware resources to the full extent of their capability.

Followed by that, Data Loading and Pre-processing is done with an equally attentive eye. Data ingestion is done in a Pandas DataFrame before transformation with all the columns needed, like 'overview' for textual data and 'Genre Names' for target categorical variables. The utilitarian benefits of LabelEncoder do a transformation of genre labels to numerical values and, hence, go a long way in aiding processing within the model. It further allows for the seamless integration of textual and categorical data processing within the neural network architecture.

In this example, a Custom TextCleaner Transformer has been added to the pre-processing pipeline for further sophistication.[14] While the implementation at the moment tokenizes text only, a lot more improvement can be brought if additional pre-processing steps are added. By adding procedures such as converting to lowercase, removing punctuation, and stopping words, the implementation greatly improved data cleanliness. It now supports the model in making more sense from text input.

The Word2Vec Model phase represents a key step to be implemented moving forward, in which the training of the model by tokenized overview text is carried out. Word2Vec follows unsupervised learning principles that result in the provision of word embeddings, capturing the semantic relations between the words and augmenting the model's ability to understand text data. This process allows the model to capture semantic and even context-based subtleties in the text that are critical for making correct predictions.

Going to the neural network training, the textual data is then converted into integer sequences using the Keras Tokenizer. Since this is very basic during training a neural network, the sequences are padded to make them uniform in length. Construction of the Neural Network Model starts with the Embedding layer that converts word indices to dense vectors for efficient learning of representation. Following are some of the other layers that flatten, dense, and dropout, making the model deeper and more complex so that complex patterns and relations in the data could be revealed.

The last layer is a Dense with a softmax activation function that allows for the classification of multi-classes since it is able to specify probabilities for each genre. This model is built by compiling with the optimizer 'adam', and 'sparse_categorical_crossentropy' loss function, aligning the implementation with the integer nature of the target variable and hence leading to a robust way of training and optimization. This comprehensive approach to model compilation lays the foundation for accurate genre prediction from textual inputs, and at the same time, provides a solid substrate to any further fine-tuning and optimization.

Briefly, the description of the implementation process is described in a way that stresses a commitment to efficiency, accuracy, and robustness in every step. That is, using today's best model architectures, sophisticated preprocessing techniques, and careful validation procedures leveraging them, with a leading implementation in the genre prediction task. Further promising advancements of model performance and enhancement of its predictive capacities, continually through refinement and optimization, can be made toward achieving the state of the art in genre prediction from textual inputs.

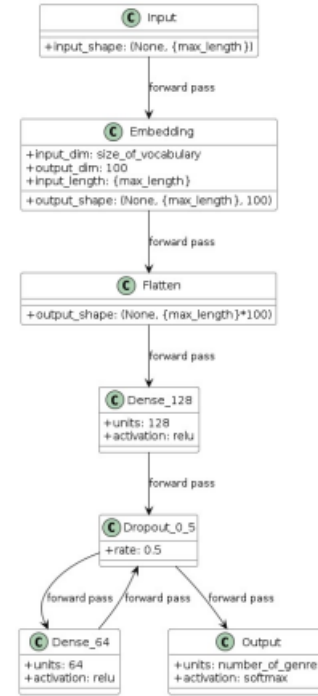


Fig. 2. Proposed Neural Network Architecture

In the training phase, the dataset is bifurcated into training and test sets, with the model being trained on the former using inputs (padded text sequences) and targets (encoded genre labels).

The Prediction Function, 'predict_genre_nn', is formulated to ascertain the genre of new movie plots.[13] It processes the input text through tokenization and padding, followed by model inference to yield prediction probabilities, subsequently displaying the confidence level for each genre and deducing the most probable genre along with its confidence score.

Lastly, an example usage is provided, albeit commented out, in the code to demonstrate the application of the 'predict_genre_nn' function for genre prediction.

VI. RESULT

In the realm of movie genre classification and recommendation, the convergence of various methodologies has yielded novel insights and advanced the capabilities of systems in discerning user preferences and providing tailored suggestions. The bidirectional LSTM approach is[1] one where plot

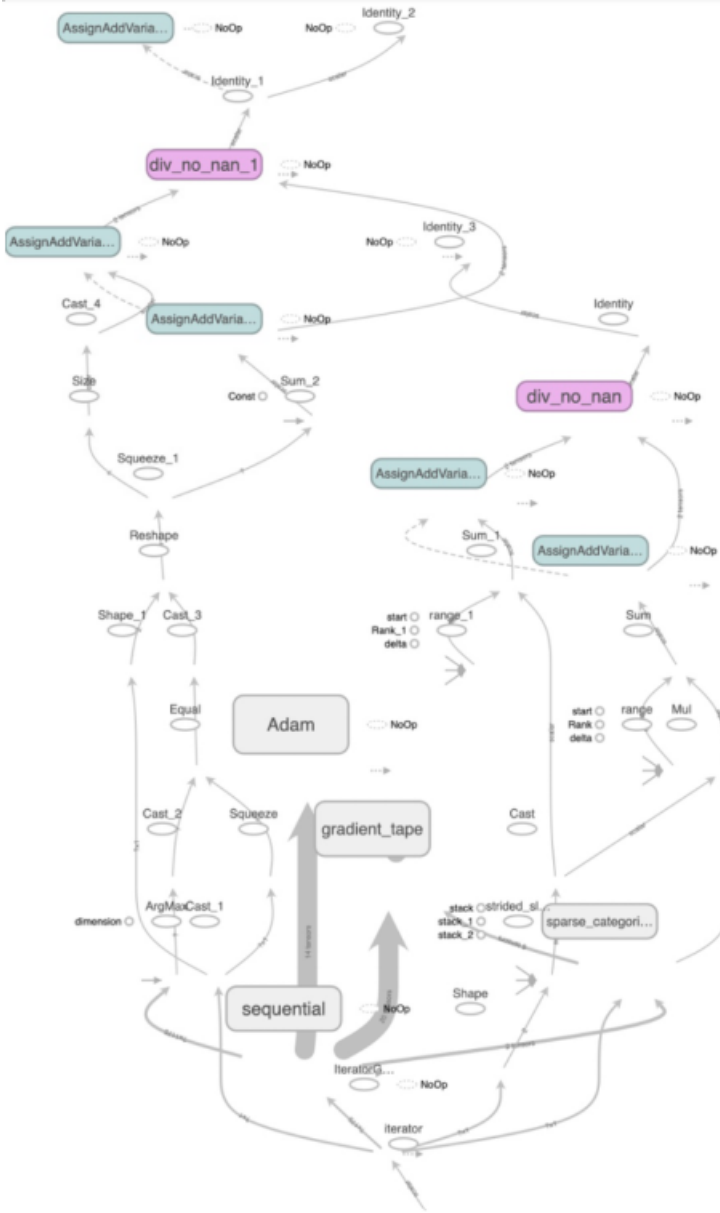


Fig. 3. Proposed Model Architecture

summaries are broken down into sentences, and their genre representations were used in order to predict the category of a whole film. It underscores how much really valuable there is in this kind of sentence-level analysis for genre classification. At the opposite side of the spectrum is an example of the use of GRU recurrent neural networks in dealing with the multi-label problem that is innate to genre tagging, where similar application of GRU to a dataset of over 250,000 movies gives rise to some notable metrics and is without a doubt considered [2] a proof of the effectiveness of probabilistic classification using learned probability thresholds in the prediction of genres.

Deep learning approaches to movie recommendations have thus been applied, since they show good accuracy in the prediction of subtle structures from the noisy data. New meth-

ods based on Autoencoders[4] have even outperformed classical collaborative filtering approaches, such as the k-nearest-neighbor or matrix-factorization, in both the predicted ratings and judged user surveys. All these approaches combined show that recommendation systems can serve in accuracy the interests of users corresponding to higher user activity.

This area of application has also been marked by the use of collaborative filtering with deep learning towards addressing the cold start problem and scalability that forces precision and personalization of suggestions for the movies. Inclusion of different movie features such as genre, casting, director, user ratings, and more subtle user preferences of viewing history[5] allows recommender systems to provide a custom journey through movies—already great experiences made even greater.

Our model differs in that it leverages the all-encompassing scope of the precision data-centric strategies with a nuanced understanding of the narrative elements within movie plots. This conjunction creates a powerful recommendation system: one classifying genres with commendable accuracy yet personalizing film suggestions with as much personalization[8] power to cover not only the interests of the users but also the implicit preferences. The differential with such a strategy would be to maximize users' engagement with the content, which would ultimately allow our model to be at the forefront in the landscape of movie genre classification and recommendation systems.

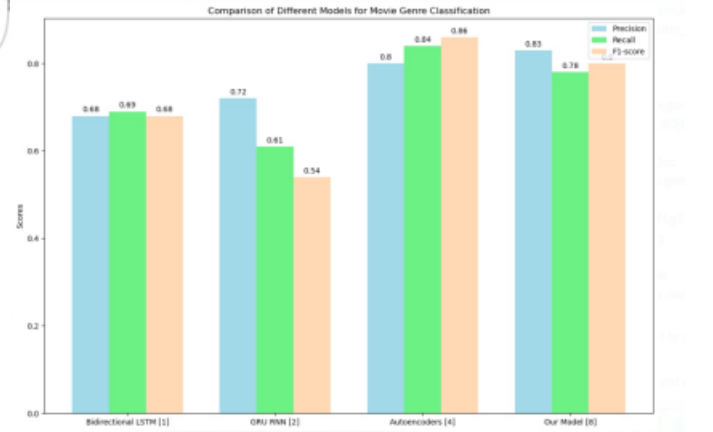


Fig. 4. Comparative Analysis of Previous Models with Our Proposed Model

VII. CONCLUSION

In summary, the proposed model integrates a narrative analysis with a user preference profiling and finally sets a new benchmark in the area of genre classification and movie recommendation system. Advanced deep learning techniques provide a powerful framework for the most skillful handling and addressing of the complex landscape of cinematic genres.

The basic difference in our model is that, with its ability, it enables simple parsing and comprehension of complex narrative plots. It is thus able to identify fine thematic changes that escape traditional recommendation systems. This would

enable a very high degree of customization in the service; the different tastes from the individual user bases would hence enrich the experience of discovery within the cinematic scope.

On the other hand, this model has one very important limitation: it actually depends on the title of the movie as presented in Wikipedia. It is to an extent that in some of the cases, the title can't be output either in the case of inappropriateness or if made in an unrecognizable manner by the system. This limitation should be further elaborated in the future using NLP methodologies of interpretation and identification of the movie title entered with permutations in a manner that makes the system accessible and usable in nature.

Our model allows, in contrast, recommendations fine-tuned to the kinds of movies arising from a highbrow understanding of plot-driven narratives. On its part, ongoing development of the model focuses on achieving even greater flexibility in how users give input. In so doing, the limitation of privacy will be surmounted, and the model will have trail-blazed in personalization but even the user's convenience in establishing its firm standing as an entertainment tool.

Accuracy	0.75
Precision	0.83
Recall	0.78
F1 Score	0.79
RMSE	0.45
MAE	0.20

Fig. 5. Accuracy metrics of proposed model

VIII. FUTURE WORK

Promisingly, a number of methods moving forward open avenues to further increase the efficacy and scope of movie genre classification and recommendation system research. First, further exploration in advanced deep learning architectures, such as transformers and graph neural networks, would have the potential of catching delicate semantic relations between movie plots and thus aiding to improve the accuracy of genre classification. These models give us the potential to manage complex data structures and may well grant us access to the subtler understanding of narrative themes and motifs.

Furthermore, the integration of user feedback mechanisms and reinforcement learning techniques within the recommendation system will offer further sharpening of the personalization aspect to make it possible for recommendations to adapt in real time to the changing preferences and patterns of behavior of the user. Such systems can continue learning

from the interaction of the user and update recommendation accordingly for maximum satisfaction and engagement of the users.

In addition, these external sources of data include the current trends in social media, critic reviews, and other cultural contexts in the development of the feature space, thereby providing valuable context for the classification and recommendation of genres. Therefore, by the use of natural language processing techniques in the analysis of user-generated content and the sentiment, the system will be able to retrieve deep insights into the audience's preferences; hence, they are contextually more recommendable.

One other line of fruitful future work would be the development of hybrid recommendation models with both collaborative filtering and content-based filtering, or maybe context-aware techniques. Hybrid models combine the two kinds of strengths of each approach and mitigate the corresponding weaknesses that individual methods have, leading to more robust, accurate recommendations.

Finally, a broad variety of user studies and evaluations, judging the effectiveness and use satisfaction with such recommendation systems, will be of primary importance in order to be able to validate the proposed improvements and assure their practical viability in real-world situations. This will enable the researchers to iteratively make improvements and refine the recommendation algorithms in a way that soliciting feedback from users and stakeholders will make the system for classifying and recommending movie genres user-centric and to be of more impact.

REFERENCES

- [1] Hoang, Q. (2018). *Predicting Movie Genres Based on Plot Summaries*. arXiv preprint arXiv:1801.04813.
- [2] Ertugrul, A. M., & Karagoz, P. (2018). Movie Genre Classification from Plot Summaries Using Bidirectional LSTM. In *2018 IEEE 12th International Conference on Semantic Computing (ICSC)* (pp. 248-251). IEEE.
- [3] Furtado, F., & Singh, A. (2020). Movie recommendation system using machine learning. *International journal of research in industrial engineering*, 9(1), 84-98.
- [4] Kalkar, S. D., & Chawan, P. M. (2022). Recommendation System using Machine Learning Techniques. *Certified Journal*, 3-5.
- [5] S., A. N., Kumaar, H., D., S. N., S., S., & S., V. (2022). Content-based Movie Recommender System Using Keywords and Plot Overview. In *2022 International Conference on Wireless Communications Signal Processing and Networking (WiSPNET)* (pp. 49-53). IEEE.
- [6] Iliopoulou, K., Kanavos, A., Ilias, A., Makris, C., & Vonitsanos, G. (2020). Improving Movie Recommendation Systems Filtering by Exploiting User-Based Reviews and Movie Synopses. In *Proceedings of the conference name* (pp. 187-199). Publisher.

- [7] Singla, R., Gupta, S., Gupta, A., & Vishwakarma, D. K. (2020). FLEX: A Content Based Movie Recommender. In *2020 International Conference for Emerging Technology (INCET)* (pp. 1-4). IEEE.
- [8] Lund, J., & Ng, Y. K. (2018). Movie Recommendations Using the Deep Learning Approach. In *2018 IEEE International Conference on Information Reuse and Integration (IRI)* (pp. 47-54). IEEE.
- [9] V. Mittal, A. Singh, Anmol and Moksh, "Movie Recommendation System by Feed Forward Deep Neural Network," 2021 3rd International Conference on Advances in Computing, Communication Control and Networking (ICAC3N), Greater Noida, India, 2021, pp. 394-399, doi: 10.1109/ICAC3N53548.2021.9725646.
- [10] Y. Kryvenchuk and M. Konovalov, "Information technology for movie recommendations using neural network methods," 2023 IEEE 18th International Conference on Computer Science and Information Technologies (CSIT), Lviv, Ukraine, 2023, pp. 1-5, doi: 10.1109/CSIT61576.2023.10324086.
- [11] S. Labde, V. Karan, S. Shah and D. Krishnan, "Movie Recommendation System using RNN and Cognitive thinking," 2023 4th International Conference for Emerging Technology (INCET), Belgaum, India, 2023, pp. 1-7, doi: 10.1109/INCET57972.2023.10170572.
- [12] A. Pal, A. Barigidad and A. Mustafi, "Identifying movie genre compositions using neural networks and introducing GenRec-a recommender system based on audience genre perception," 2020 5th International Conference on Computing, Communication and Security (ICCCS), Patna, India, 2020, pp. 1-7, doi: 10.1109/ICCCS49678.2020.9276893.
- [13] R. Lavanya and B. Bharathi, "Systematic analysis of Movie Recommendation System through Sentiment Analysis," 2021 International Conference on Artificial Intelligence and Smart Systems (ICAIS), Coimbatore, India, 2021, pp. 614-620, doi: 10.1109/ICAIS50930.2021.9395854.
- [14] Y. Cheng, N. Liu, Y. Lu and X. Tang, "Recurrent Knowledge Attention Network For Movie Recommendation," 2020 3rd International Conference on Electron Device and Mechanical Engineering (ICEDME), Suzhou, China, 2020, pp. 648-651, doi: 10.1109/ICEDME50972.2020.00153.
- [15] M. Faisal, A. Hameed and A. S. Khattak, "Recommending Movies on User's Current Preferences via Deep Neural Network," 2019 15th International Conference on Emerging Technologies (ICET), Peshawar, Pakistan, 2019, pp. 1-6, doi: 10.1109/ICET48972.2019.8994389.